

Juliana Porto

March 30, 2024

DS503 Big Data Management

Project 3

Question	Status	Comment
Q1	Fully Working	Step 1: Creating Dataset I used the same transaction dataset that I created for Project 1. Source Code: Project3/Q1_SparkSQL_TransactionDataProcessing/Main.java Transactions Dataset (T): Project3/Q1_SparkSQL_TransactionDataProcessing/Transactions.txt
		Step 2: Spark Workflow I completed the workflow in PySpark using Google Collab. I found this easiest as I had multiple issues with Scala and my IDE. <i>Plan:</i> <ol style="list-style-type: none">1. Setup Spark2. Import transaction data3. Add column headers (check schema & dataframe)4. Start with T5. Create T1<ol style="list-style-type: none">a. Start with Tb. Filter<ol style="list-style-type: none">i. TransTotal >= 2006. Create T2<ol style="list-style-type: none">a. Start with T1b. GroupBy TransNumItemsc. Aggregate<ol style="list-style-type: none">i. Sum of TransTotalii. Avg of TransTotaliii. Min of TransTotaliv. Max of TransTotald. OrderBy7. Report T2 to client using show8. Create T3<ol style="list-style-type: none">a. Start with T1

		<ul style="list-style-type: none"> b. GroupBy CustID c. Aggregate
		<ul style="list-style-type: none"> <ul style="list-style-type: none"> i. Count of CustID d. Name it as NumTransT3 for later use e. Select <ul style="list-style-type: none"> i. CustID ii. NumTransT3 9. Report T3 to client using show 10. Create T4 <ul style="list-style-type: none"> a. Start with T b. Filter <ul style="list-style-type: none"> i. TransTotal >= 600 11. Create T5 <ul style="list-style-type: none"> a. Start with T4 b. GroupBy CustID c. Aggregate <ul style="list-style-type: none"> i. Count of CustID d. Name it as NumTransT5 for later use e. Select <ul style="list-style-type: none"> i. CustID ii. NumTransT5 12. Report T5 to client using show 13. Create T6 <ul style="list-style-type: none"> a. Create T5 with T3 join <ul style="list-style-type: none"> i. Inner join on CustID b. Start with T5 with T3 join <ul style="list-style-type: none"> i. Filter <ul style="list-style-type: none"> 1. NumTransT5 * 5 < NumTransT3 ii. Select <ul style="list-style-type: none"> 1. CustID 14. Report T6 to client using show <p><i>Note: My Transaction Data did not end up having any customers in T6, so I tested it with a separate smaller test dataset and included the screenshot (T6_verify) of that output to show that it does indeed work when there is data that fits that condition.</i></p>

		<p>Source Code: Project3/Q1_SparkSQL_TransactionDateProcessing/ DS503_Project3_Q1.ipynb Project3/Q1_SparkSQL_TransactionDateProcessing/ ds503_project3_q1.py <i>Note:</i> Can also be accessed at: https://colab.research.google.com/drive/1Xhfhel3j3XklwbhFnQLZrGdkta6QYttO?usp=sharing</p>
		<p>Output with Transaction Dataset: Project3/Q1_SparkSQL_TransactionDateProcessing/ Output_Screenshots/</p> <p>Additional Notes:</p> <ul style="list-style-type: none"> Used SparkSQL & DataFrames to write the workflow Ran into multiple Scala issues because I was unfamiliar with the language and my IDE would not cooperate after many attempts – my solution was to use PySpark as I am familiar with python. Overall, I found that using PySpark instead of Scala allowed me to focus more on how to use SparkSQL as opposed to having to figure out the details of Scala.
Q2	Fully Working	<p>Step 1: Dataset Used the point dataset created in Project 2. I also created a smaller dataset with a small number of points in specific areas for testing purposes. <i>Assumptions:</i></p> <ul style="list-style-type: none"> All values are integers <p>Points Dataset: Project3/Q2_SparkRDDs_GridCellsofHighRelativeDensityIndex/ Points.txt</p> <p>Test Point Datasets: .../Test_Points.txt</p> <p>Step 2 & 3: Report the TOP 50 grid cells w.r.t Relative-Density Index & Neighbors of the TOP 50 grids <i>Plan:</i></p> <ol style="list-style-type: none"> 1. Import points dataset as data frame and set schema 2. Add column headers 3. Assign point's current grid

		<div><div><div>4. Get count of points in each grid</div><div>5. Assign neighbors & check if valid</div><div>6. Get count of neighbors</div><div>7. For each grid, calculate density</div><div>8. Display final outputs</div><div><div>a. Top 50 grids with highest relative-density index</div><div>b. Neighbors of top 50 grids</div></div></div><div><div>Source Code:</div><div>Test Dataset:</div><div>Project3/Q2_SparkRDDs_GridCellsofHighRelativeDensityIndex/DS503_Project3_Q2_Test.ipynb</div></div></div>																
		<div><div>Project3/Q2_SparkRDDs_GridCellsofHighRelativeDensityIndex/ds503_project3_q2_test.py</div><div>Note: Can also be accessed at:</div><div>https://colab.research.google.com/drive/1Wrx8eRfWII_X_KAyjBT4ivy_2f17AUbV?usp=sharing</div><div>Points Dataset:</div><div>Project3/Q2_SparkRDDs_GridCellsofHighRelativeDensityIndex/ DS503_Project3_Q2.ipynb</div><div>Project3/Q2_SparkRDDs_GridCellsofHighRelativeDensityIndex/ds503_project3_q2.py</div><div>Note: Can also be accessed at:</div><div>https://colab.research.google.com/drive/1rH9bRHKEoD4703EMgYOZAP5uem95zacP?usp=sharing</div><div>Output with Test Dataset:</div><div>Project3/Q2_SparkRDDs_GridCellsofHighRelativeDensityIndex/Outputs/Test_Top50Grid</div><div>Output with Points Dataset:</div><div>Project3/Q2_SparkRDDs_GridCellsofHighRelativeDensityIndex/Outputs/Points_Top50Grid</div></div> <div><div>Additional Notes:</div><div><div>• Grid for reference:</div><table><tr><td>249,501</td><td>249,502</td><td>...</td><td>250,000</td></tr><tr><td>...</td><td>...</td><td>...</td><td>...</td></tr><tr><td>501</td><td>502</td><td>...</td><td>1000</td></tr><tr><td>1</td><td>2</td><td>...</td><td>500</td></tr></table></div></div>	249,501	249,502	...	250,000	501	502	...	1000	1	2	...	500
249,501	249,502	...	250,000															
...															
501	502	...	1000															
1	2	...	500															

		<ul style="list-style-type: none">• If a neighbor does not exist, its grid value is labeled as -1 and is not considered for the average calculation, but does appear in the data frame• Once again used PySpark in google collab as I was unable to get Scala working with my IDE• Final neighbor output from points dataset took a very long time to run and I did not get to include the final output screenshot because of this, but you should be able to see it through the google collab link
--	--	--